

HLB: Benchmarking LLMs’ Humanlikeness in Language Use

Xufeng Duan¹ Bei Xiao¹ Xuemei Tang¹ Zhenguang G. Cai^{1,2}

¹Department of Linguistics and Modern Languages, The Chinese University of Hong Kong

²Brain and Mind Institute, The Chinese University of Hong Kong

xufeng.duan@link.cuhk.edu.hk

Abstract

As synthetic data becomes increasingly prevalent in training language models, particularly through generated dialogue, concerns have emerged that these models may deviate from authentic human language patterns, potentially losing the richness and creativity inherent in human communication. This highlights the critical need to assess the humanlikeness of language models in real-world language use. In this paper, we present a comprehensive humanlikeness benchmark (HLB) evaluating 20 large language models (LLMs) using 10 psycholinguistic experiments designed to probe core linguistic aspects, including sound, word, syntax, semantics, and discourse (see [this link](#)). To anchor these comparisons, we collected responses from over 2,000 human participants and compared them to outputs from the LLMs in these experiments.

For rigorous evaluation, we developed a coding algorithm that accurately identified language use patterns, enabling the extraction of response distributions for each task. By comparing the response distributions between human participants and LLMs, we quantified humanlikeness through distributional similarity. Our results reveal fine-grained differences in how well LLMs replicate human responses across various linguistic levels. Importantly, we found that improvements in other performance metrics did not necessarily lead to greater humanlikeness, and in some cases, even resulted in a decline. By introducing psycholinguistic methods to model evaluation, this benchmark offers the first framework for systematically assessing the humanlikeness of LLMs in language use.

1 Introduction

In recent years, large language models (LLMs) have made significant advancements. Models like OpenAI’s GPT series and Meta’s Llama family can

generate human-like text, engage in coherent dialogues, and answer complex questions, often producing responses that are indistinguishable from those of humans in certain evaluations (Tsubota and Kano, 2024). Cai et al. (2024) conducted a systematic evaluation of human-like language use in models such as ChatGPT and Vicuna, demonstrating that LLMs closely replicate human language patterns in many aspects. However, despite these successes, questions remain about how accurately these models capture the deeper, nuanced patterns of human language use. In other words, the full extent of their similarity to human behavior remains unclear.

The importance of evaluating humanlikeness in language use is further underscored by the increasing reliance on synthetic data for model training, particularly in dialogue models. While synthetic data generation facilitates efficient scaling of model training, it raises concerns about models diverging from real-world human language patterns (del Rio-Chanona et al., 2024). Studies have shown that synthetic data can degrade model performance after retraining (Shumailov et al., 2024). This makes it imperative to assess the humanlikeness of LLMs rigorously across various aspects of language use, to ensure that models do not lose the diversity and richness of human language data.

To address this challenge, we introduce a psycholinguistic benchmark designed to provide a systematic and comprehensive evaluation of how closely LLMs align with human linguistic behavior.

Although numerous benchmarks and leaderboards have been developed to assess the performance of LLMs on downstream NLP tasks, they often fail to capture the finer, human-like qualities of language use. Current NLP benchmarks typically focus on task-based accuracy or performance (Lewkowycz et al., 2022; Zhou et al., 2023; Peng et al., 2024; Hendrycks et al., 2021; Zellers et al.,

2019), overlooking the broader psycholinguistic dimensions that characterize how humans process and produce language. Furthermore, few studies have systematically compared the language use of LLMs and human participants across multiple linguistic levels. This gap highlights the need for a new benchmark that can robustly measure the extent to which LLMs replicate human language behavior in real-world, diverse linguistic contexts.

In this paper, we address this gap by presenting a psycholinguistic benchmark study that evaluates the humanlikeness of 20 LLMs. Our benchmark consists of 10 representative psycholinguistic experiments, adapted from Cai et al. (2024), which cover five core linguistic aspects: sound, word, syntax, semantics, and discourse, with two experiments dedicated to each aspect (see 1). We collected approximately 50 to 100 responses per item from over 2,000 human participants. Additionally, we gathered 100 responses per item from each of the 20 LLMs, including well-known models such as GPT-4o, GPT-3.5, Llama 2, Llama 3, Llama 3.1, and other state-of-the-art models (see Table 1). To quantify humanlikeness, we developed an auto-coding algorithm that efficiently and reliably extracts language use patterns from responses. The humanlikeness metric was then calculated based on the similarity between the response distributions of humans and LLMs, using a comparison of their probability distributions.

Our findings reveal significant, nuanced differences in how LLMs perform across various linguistic aspects, offering a new benchmark for evaluating the humanlikeness of LLMs in natural language use. This benchmark introduces psycholinguistic methods to model evaluation and provides the first framework for systematically assessing the humanlikeness of LLMs in language use.

2 Related Work

Recent advances in LLMs have led to the development of various benchmarks designed to evaluate their linguistic capabilities. Standard benchmarks like GLUE (Wang et al., 2018) and SuperGLUE (Wang et al., 2019) assess models across a range of natural language processing (NLP) tasks, including sentence classification, textual entailment, and question answering. However, these benchmarks primarily focus on task-based accuracy and often overlook the more intricate aspects of human-like language processing. While these evaluations

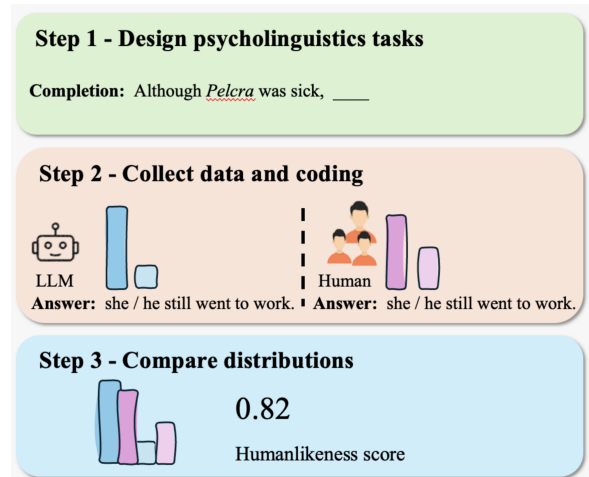


Figure 1: The benchmark framework. The example prompt is taken from the sound-gender association task, where humans can infer the gender of a novel name (e.g., *Pelcra* for Female; *Pelcrad* for Male) based on phonology.

provide valuable insights into model performance, they do not fully capture the extent to which LLMs comprehend and generate language in a humanlike manner. As Manning et al. (2020) note, LLMs are powerful statistical models that can identify patterns in vast datasets, but these benchmarks do not adequately test how well models replicate human patterns of language use due to the interplay of complex cognitive biases.

2.1 Psychological Experimentation on LLMs

A growing body of research has begun applying classical psychological experiments to evaluate LLMs in more domain-specific and cognitively demanding tasks. For example, Binz and Schulz (2023) and Dasgupta et al. (2023) used well-known psychological paradigms, such as the Linda problem and the Wason selection task, to probe LLMs' abilities in judgment and decision-making. Similarly, Sap et al. (2023) and Trott et al. (2023) explored whether LLMs exhibit theory of mind, a key component of human social cognition, while Miotto et al. (2022) and Karra et al. (2023) examined LLMs' personality traits. In the domain of behavioral economics, Horton (2023) conducted experiments with GPT-3 to explore its decision-making processes. These studies suggest that LLMs can be treated as cognitive agents in psychological experiments, providing insights into how LLMs align with humans in reasoning, behavior, and decision-making. Moreover, they help shed light on the underlying mechanisms of LLMs,

Exp Level	Task
E1	sound-shape association (Köhler, 1967)
E2	Sound sound-gender association (Cassidy et al., 1999)
E3	Word word length and predictivity (Mahowald et al., 2013)
E4	Word word meaning priming (Rodd et al., 2013)
E5	structural priming (Pickering and Branigan, 1998)
E6	Syntax syntactic ambiguity resolution (Altmann and Steedman, 1988)
E7	Meaning implausible sentence interpretation (Gibson et al., 2013)
E8	semantic illusion (Erickson and Mattson, 1981)
E9	Discourse implicit causality (Garvey and Caramazza, 1974)
E10	Discourse drawing inferences (Singer and Spear, 2015)

Table 1: The experiments in this benchmark.

as seen in the work of Huang and Chang (2023) and Qiao et al. (2023), who analyzed reasoning patterns in LLMs. Hagendorff (2023) further provided a comprehensive review of LLM performance in psychological tests, showing that while LLMs demonstrate sophisticated behaviors, they often diverge from human cognition. These divergences highlight the need for more robust frameworks to understand the limitations of LLMs in mimicking human thought processes.

2.2 Psycholinguistic Experimentation on LLMs

Psycholinguistic approaches offer a deeper analysis by testing LLMs on how well they replicate the cognitive processes underlying human language processing. Ettinger (2020) and Futrell (2019) have subjected models like BERT to psycholinguistic tasks such as syntactic ambiguity resolution and structural priming, revealing both the strengths and limitations of LLMs in replicating human language processing. Michaelov and Bergen (2023) used

structural priming tasks to investigate how LLMs internalize syntactic structures, while Huang et al. (2024) examined LLMs’ ability to resolve syntactic ambiguity. Qiu et al. (2023) explored how well LLMs handle pragmatic reasoning. These studies demonstrate that LLMs can, to some extent, mimic humanlike behavior in controlled experiments. However, divergences in processing reveal the distinctions between machine learning models and humans. A recent review by Demszky et al. (2023) emphasized the need for benchmarks that incorporate psychological paradigms to evaluate LLMs. The authors argue that by applying psycholinguistic methods, researchers can better understand how closely LLMs approximate human cognition and where they fall short. Despite extensive research on LLMs’ performance across various tasks, there is still no benchmark that includes human language processing data to reveal the extent to which LLMs resemble humans, particularly in language use. This paper addresses that gap by adapting 10 psycholinguistic experiments from Cai et al. (2024) to evaluate how closely LLMs align with human language behavior, covering phenomena ranging from sound symbolism to discourse comprehension.

3 Methodology

3.1 Human Experiments

Experimental Design The human experiments were conducted using Qualtrics, an online survey platform (Qualtrics, 2024). The study included ten psycholinguistic tasks that spanned various linguistic levels, from sound, word, syntax, and meaning to discourse comprehension, with two experiments for each level (see Appendix A for details). We exposed a participant to only one trial on each experiment, with a total of 10 trials across all the experiments. This setup minimized trial-level effects and facilitated direct comparisons with LLMs, which were tested under similar conditions (presenting instructions and stimuli in a single prompt) to avoid context effects within individual conversations.

Procedure After providing consent, participants completed the ten psycholinguistic tasks (presented in a random order); four attention checks were randomly interspersed among the trials to later identify participants for random responding. Each experimental task began with an instructional screen, some of which included examples to clarify task

requirements. The examples were carefully designed to differ from the experimental stimuli to prevent potential priming effects. For instance, in a sentence-completion task, an illustrative example that did not resemble the experimental stimuli and did not induce target words for any stimuli was used. The priming tasks (which included pairs of priming and target stimuli) were spread across multiple pages to avoid strategic responses in case participants realise the relation between the prime and the target. The overall experimental procedure was streamlined for clarity and efficiency, with each session lasting approximately 8 to 10 minutes (mean = 8.336, $SD = 4.171$).

Participants Participants were recruited from the crowd-sourcing platform Quattrics and restricted to native English speakers residing in the UK and US, according to their registration on Prolific. They were required to use a desktop computer to complete the tasks. Among the 2,205 participants taking part in the experiments, 290 were excluded for not well adhering to the experimental instructions, including completing the study too quickly, showing low effort, or not finishing the experiment, according to the Qualtrics system. The remaining 1,915 participants were further checked for language nativeness and their accuracy with attention checks. After a thorough screening process—excluding those who were not native speakers, failed attention checks, or exhibited irregularities such as excessively short completion times or multiple participation attempts—the final valid sample consisted of 1,905 participants. The sample was composed of participants as follows: female ($n = 1,051$), male ($n = 838$), preferred not to disclose ($n = 16$), with an average age of 44.8 years (range: 18 to 89 years). Educational levels included: no formal education ($n = 2$), elementary school ($n = 12$), high school ($n = 672$), bachelor’s degree ($n = 862$), and master’s degree ($n = 357$). This sample of participants resulted in each item being tested in a minimum average of 24 trials (e.g., Word Length and Predictability) and up to an average of 96 trials (e.g., Sound-Shape Association Task).

3.2 LLM Experiments

Experimental Design To compare human responses with those generated by LLMs, we employed the same 10 psycholinguistic tasks designed for human participants. 20 LLMs (See Table 2) were selected for evaluation, including models from prominent families like OpenAI’s GPT series

(GPT-4o, GPT-3.5), Meta’s Llama series (Llama 2, Llama 3, Llama 3.1) and Mistral series (OpenAI et al., 2024; Touvron et al., 2023; AI, 2024). Each model provided 100 responses per item in each experiment, ensuring that the response data was comparable to the human data. Similar to the human experimental design, LLMs followed a one-trial-per-run paradigm, ensuring that responses were generated independently for each item to prevent context effects. The input format for the LLMs closely mirrored the instructions provided to human participants. Careful modification of human prompts was performed to ensure that task instructions were clear and interpretable by LLMs. This allowed for a direct comparison between human and LLM performance on the same tasks under identical conditions.

Response Collection Procedure This closely mirror that in the human experiments. Each LLM was presented with the task instructions and the stimulus combined into a single prompt. We collected 100 responses (across different conditions) for each item in an experiment in order to ensure a sufficiently large dataset for robust analysis of the response distributions. For OpenAI models, responses were obtained through the OpenAI API, while models hosted on Hugging Face were accessed using the Hugging Face Inference API. All requests to the models were made using their default parameters to encourage variability in responses. The collected responses were stored and processed for subsequent coding and analysis.

3.3 Response Coding

Development and Validation We employed an auto-coding algorithm across 10 experiments to assess agreement between human annotations and machine-generated labels. This algorithm utilized spaCy’s *en_core_web_trf-3.7.3* model for syntactic parsing (e.g., structural priming and syntactic ambiguity resolution tasks) and regular expressions to detect answer patterns in others. Across 20,953 trials of human response data, we computed Cohen’s Kappa (κ), a measure that corrects for chance agreement between the results from manually coding and auto-coding algorithm, defined as:

$$K = \frac{P_o - P_e}{1 - P_e} \quad (1)$$

where P_o is the observed agreement, and P_e is the expected agreement by chance.

The Kappa score was $\kappa = 0.993$, indicating near-perfect agreement ($z = 451, p < 0.001$). This demonstrates the high accuracy of the auto-coding algorithm in replicating human annotations.

3.4 Humanlikeness Scoring

To quantify the humanlikeness of LLM responses, we used Jensen-Shannon (JS) divergence to compare the response distributions between human participants and LLMs. JS divergence, a symmetric measure of similarity between two probability distributions, is ideal for assessing how closely LLM responses mirror human behavior across linguistic levels. For each task, the auto-coding algorithm generated response distributions for both humans and LLMs. We computed **humanlikeness score (HS)** for each item as:

$$\begin{aligned} HS_{\text{item}} &= 1 - JS(P, Q) \\ &= 1 - \frac{1}{2} [KL(P \parallel M) + KL(Q \parallel M)] \end{aligned} \quad (2)$$

where P and Q are the human and LLM response distributions, and M is their average. For each experiment, we average the scores across all items. The overall humanlikeness score across all experiments is then computed as:

$$HS_{\text{Overall}} = \frac{1}{m} \sum_{j=1}^m \left(\frac{1}{n_j} \sum_{i=1}^{n_j} \left(1 - \frac{1}{2} [KL(P_i \parallel M_i) + KL(Q_i \parallel M_i)] \right) \right) \quad (3)$$

4 Result

The overall humanlikeness scores revealed notable variations in how well LLMs emulated human language use across the 10 psycholinguistic experiments. Here, we performed a concise analysis to explore the data.

OpenAI’s models, including GPT-3.5-turbo, GPT-4o-mini, and GPT-4o, exhibited relatively stable performance across tasks, maintaining consistent humanlikeness scores. In contrast, the Llama family of models showed an overall increase in humanlikeness scores, with Meta-Llama-3.1-70B-Instruct achieving the highest performance among all Llama models. On the other hand, the Mistral family of models showed a slight decrease in humanlikeness, with Mistral-7B-Instruct-v0.3 scoring lower than its predecessors, indicating less alignment with human language use.

4.1 Comparative Analysis

Statistical comparisons between model families (three models selected per model family) revealed significant differences in performance. Notably, Llama models significantly outperformed Mistral models in humanlikeness ($t = 10.44, p < .001$) highlighting the substantial gap between these two families. Furthermore, Llama models also outperformed OpenAI models ($t = 3.13, p = .002$) although this difference was less pronounced compared to the Llama vs. Mistral comparison.

Within the Llama family, the transition from Meta-Llama-3-70B-Instruct to Meta-Llama-3.1-70B-Instruct showed a significant increase in humanlikeness ($t = -4.85, p < .001$), indicating improvements in model performance. In contrast, no significant differences were observed between GPT-3.5-turbo and GPT-4o ($t = -0.93, p = 0.352$), suggesting that OpenAI’s models performed consistently across experiments. Interestingly, within the Mistral family, Mistral-7B-Instruct-v0.3 showed a significant decrease in humanlikeness compared to Mistral-7B-Instruct-v0.1 ($t = 5.45, p < .001$).

These results underscore the varying abilities of different model families to approximate human language patterns, with Llama models demonstrating superior performance overall.

4.2 Case analysis

An in-depth analysis of individual experiments further highlights how LLMs’ performance varies in replicating human-like responses. Experiment 4, which tested word meaning priming, emerged as the most non-humanlike among the tasks, with substantial differences between human and LLMs’ responses ($t = -116.32, p < .001$). In this experiment, we assessed whether humans and models tend to access, when reading an ambiguous word such as *post*, the meaning previously used in the prime of an ambiguous word. Human participants exhibited a modest priming effect, with 20% associating *post* with its job-related meaning after the word-meaning prime and 18% after the synonym prime. In contrast, the Llama-3.1-70B model demonstrated a significantly higher priming effect, with 52% responding to the word-meaning prime and 38% to the synonym prime, revealing a stark divergence from human patterns. This case study emphasizes the challenges LLMs face in aligning their semantic associations with human interpretations, particularly when processing ambiguous or

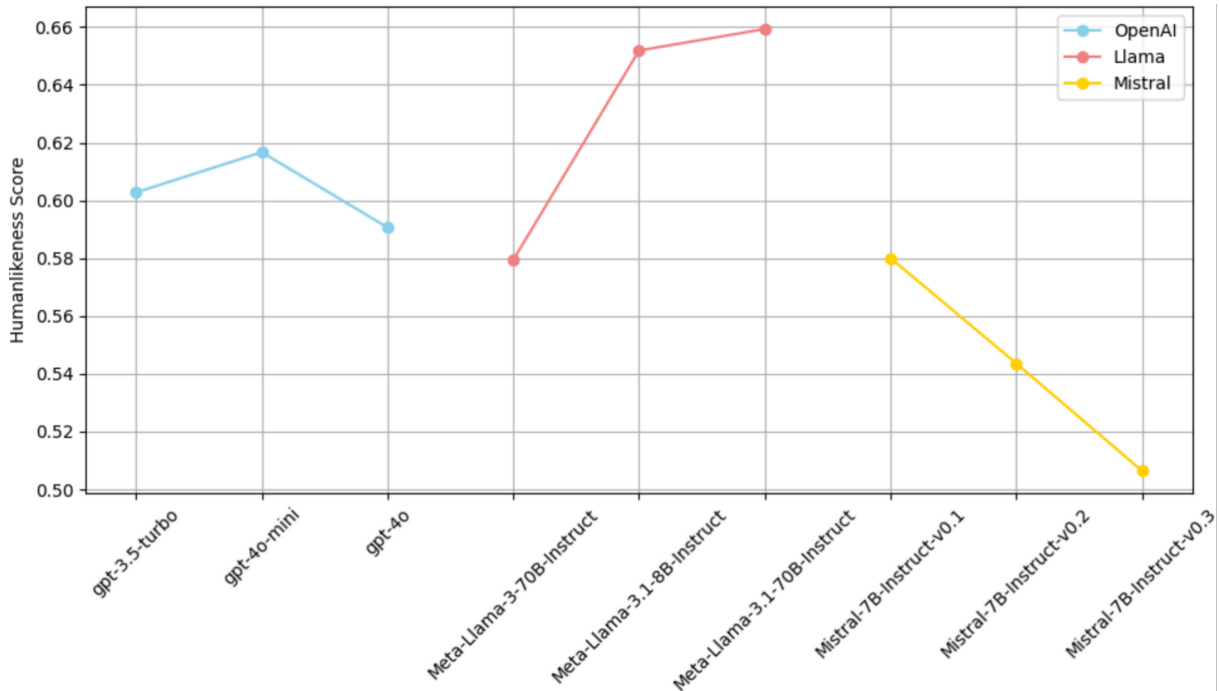


Figure 2: Humanlikeness scores of three LLM families

polysemous words.

5 Discussion

The results of this benchmark study highlight notable differences in how LLMs approximate human language use across various linguistic levels. The Llama family of models, particularly Meta-Llama-3.1-70B-Instruct, consistently outperformed both the OpenAI and Mistral models in terms of humanlikeness score. This finding suggests that recent advancements in the Llama models have led to more humanlike language behaviors, especially in terms of semantic and discourse processing. The OpenAI models, including GPT-4o and GPT-3.5-turbo, showed relatively stable performance across tasks, with no significant differences between the models. This stability may reflect a plateau in the improvement of humanlikeness in these models, as compared to the more recent gains observed in the Llama family. On the other hand, the Mistral models demonstrated a decrease in humanlikeness scores, particularly in the transition to Mistral-7B-Instruct-v0.3. This suggests that certain training methods and data quality in Mistral may have reduced their alignment with human language patterns. One of the key insights from this study is that models differ not only in their overall humanlikeness scores but also in how they handle specific linguistic phenomena. For instance, in Experiment

4 (word meaning priming), we observed a significant divergence in responses between humans and LLMs, with the latter showing a much larger priming effect. This over-priming suggests that while LLMs may excel in certain aspects of language generation, they often lack the subtle flexibility that humans display when processing ambiguous or context-dependent language. A major strength of this study is its use of psycholinguistic experiments to evaluate LLMs, which goes beyond traditional NLP benchmarks that focus on task accuracy. By systematically probing various linguistic levels—sound, word, syntax, semantics, and discourse—this benchmark provides a more comprehensive understanding of how LLMs process and generate language.

6 Conclusion

In this paper, we introduced a novel benchmark for evaluating the humanlikeness of LLMs in language use based on psycholinguistic experiments. Our study evaluated 20 LLMs, including OpenAI’s GPT family, Meta’s Llama family, the Mistral family and others, across 10 experiments that spanned key linguistic aspects such as sound, word, syntax, semantics, and discourse. Using responses from over 2,000 human participants as a baseline, the results revealed significant differences in model performance, with Llama models consistently out-

Experiment	Overall	Sound		Word		Meaning		Syntax		Discourse	
		E1	E2	E3	E4	E5	E6	E7	E8	E9	E10
Meta-Llama-3.1-70B-Instruct	66.50	89	62	61	6	81	77	80	67	80	63
Meta-Llama-3.1-8B-Instruct	65.89	73	65	60	12	84	78	79	74	79	56
Phi-3-mini-4k-instruct	64.61	61	68	59	19	89	71	76	48	80	76
Mistral-Nemo-Instruct-2407	63.69	74	63	56	10	84	77	52	75	79	68
Llama-2-13b-chat-hf	63.15	57	57	51	25	75	67	74	72	76	79
Mistral-7B-Instruct-v0.1	62.77	73	70	62	24	87	43	69	36	79	84
CodeLlama-34b-Instruct-hf	62.18	79	64	60	23	82	53	58	63	79	61
c4ai-command-r-plus	60.77	79	60	63	8	72	72	66	59	78	50
Meta-Llama-3-8B-Instruct	60.65	69	53	57	14	79	78	59	66	77	54
starchat2-15b-v0.1	60.57	58	70	58	25	87	73	62	36	75	62
gpt-4o	58.58	60	63	68	2	71	77	47	61	75	62
gpt-3.5-turbo	58.32	55	61	66	3	76	76	71	47	76	50
Yi-1.5-34B-Chat	58.20	67	54	55	13	72	70	61	65	78	48
Llama-2-7b-chat-hf	57.47	74	61	58	22	67	69	50	60	75	39
zephyr-7b-alpha	56.96	57	62	47	23	85	29	44	73	76	75
Meta-Llama-3-70B-Instruct	56.73	60	61	55	4	71	75	57	59	75	50
gpt-4o-mini	56.21	56	58	62	3	70	75	46	58	75	57
Mistral-8x7B-Instruct-v0.1	52.80	60	53	48	23	71	46	43	59	73	52
Mistral-7B-Instruct-v0.3	52.45	53	58	47	25	75	38	49	59	73	47
Mistral-7B-Instruct-v0.2	50.18	13	58	54	14	72	61	46	64	71	49
zephyr-7b-beta	47.85	28	53	48	26	71	7	38	73	75	60

Table 2: The humanlikeness score for models in different experiments.

performing both OpenAI and Mistral models in terms of language use humanlikeness. These findings underscore the potential of psycholinguistic benchmarks to capture aspects of language that are often missed by traditional NLP evaluations.

This benchmark provides a framework for future research on LLMs, offering a more meaningful and comprehensive way to evaluate their performance in real-world language use. It also highlights areas where current LLMs diverge from human language patterns, particularly in tasks involving semantic priming and ambiguity resolution. By identifying these gaps, this study offers critical insights for the next generation of LLM development, paving the way for models that more closely mirror the intricacies of human communication.

7 Limitation

However, there are several limitations to this study. First, while the benchmark covers a wide range of linguistic tasks, it may not encompass the full complexity of human language use. Some linguistic phenomena, such as pragmatic reasoning, were not explored in this study. Second, we did not manipulate models’ parameters, particularly the temperature or top k, to control the diversity of the generated responses. While using default parameters, particularly temperature, may seem limiting, this choice ensures that we evaluate models in their most typical and practical configurations. Default settings reflect how these models are commonly used in real-world applications, offering a fair and standardized comparison. Tuning parameters like temperature could introduce bias and variability across models, making it difficult to ensure con-

sistent evaluation. By using default settings, we eliminate these concerns, allowing for a more reliable assessment of humanlikeness. Finally, while the study includes a large sample of human participants, the specific demographic characteristics (e.g., native English speakers from the UK and US) may not fully represent global language use patterns. Compared to previous benchmarks that focus on task-based performance, this study offers a more in-depth analysis of language models' alignment with human linguistic behavior. Similar studies, such as Ettinger (2020), have used psycholinguistic principles to probe LLMs, but our study stands out by incorporating a broader range of linguistic levels and by using a large-scale dataset of human responses for direct comparison. The significant differences found between model families, such as the higher humanlikeness of Llama models, provide valuable insights for the ongoing development and fine-tuning of LLMs.

References

- Mistral AI. 2024. Mistral-7b-instruct-v0.3: An advanced instruction-based language model. Hugging Face Model Card. Released on May 22, 2024. Available at: <https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.3>.
- Gerry Altmann and Mark Steedman. 1988. Interaction with context during human sentence processing. *Cognition*, 30(3):191–238.
- Marcel Binz and Eric Schulz. 2023. [Using cognitive psychology to understand GPT-3](#). *Proceedings of the National Academy of Sciences*, 120(6):e2218523120. Publisher: Proceedings of the National Academy of Sciences.
- Zhenguang G. Cai, Xufeng Duan, David A. Haslett, Shuqi Wang, and Martin J. Pickering. 2024. [Do large language models resemble humans in language use?](#) *arXiv preprint*. ArXiv:2303.08014 [cs].
- Kimberly Wright Cassidy, Michael H Kelly, and Lee'at J Sharoni. 1999. Inferring gender from name phonology. *Journal of Experimental Psychology: General*, 128(3):362.
- Ishita Dasgupta, Andrew K. Lampinen, Stephanie C. Y. Chan, Hannah R. Sheahan, Antonia Creswell, Dharshan Kumaran, James L. McClelland, and Felix Hill. 2023. [Language models show human-like content effects on reasoning tasks](#). *arXiv preprint*. ArXiv:2207.07051 [cs].
- R Maria del Rio-Chanona, Nadzeya Laurentsyeva, and Johannes Wachs. 2024. [Large language models reduce public knowledge sharing on online Q&A platforms](#). *PNAS Nexus*, page pgae400.
- Dorottya Demszky, Diyi Yang, David S. Yeager, Christopher J. Bryan, Margaret Clapper, Susannah Chandhok, Johannes C. Eichstaedt, Cameron Hecht, Jeremy Jamieson, Meghann Johnson, Michaela Jones, Danielle Krettek-Cobb, Leslie Lai, Nirel Jones Mitchell, Desmond C. Ong, Carol S. Dweck, James J. Gross, and James W. Pennebaker. 2023. [Using large language models in psychology](#). *Nature Reviews Psychology*, 2(11):688–701. Publisher: Nature Publishing Group.
- TD Erickson and ME Mattson. 1981. From words to meaning: A semantic illusion. *J verbal learn verbal behav* 20 (5): 540–551.
- Allyson Ettinger. 2020. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.
- R Futrell. 2019. Neural language models as psycholinguistic subjects: Representations of syntactic state. *arXiv preprint arXiv:1903.03260*.
- Catherine Garvey and Alfonso Caramazza. 1974. Implicit causality in verbs. *Linguistic inquiry*, 5(3):459–464.
- Edward Gibson, Leon Bergen, and Steven T Piantadosi. 2013. Rational integration of noisy evidence and prior semantic expectations in sentence interpretation. *Proceedings of the National Academy of Sciences*, 110(20):8051–8056.
- Thilo Hagendorff. 2023. [Machine Psychology: Investigating Emergent Capabilities and Behavior in Large Language Models Using Psychological Methods](#). *arXiv preprint*. ArXiv:2303.13988 [cs].
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring Massive Multitask Language Understanding](#). *arXiv preprint*. ArXiv:2009.03300 [cs].
- Jie Huang and Kevin Chen-Chuan Chang. 2023. [Towards Reasoning in Large Language Models: A Survey](#). *arXiv preprint*. ArXiv:2212.10403 [cs].
- Kuan-Jung Huang, Suhas Arehalli, Mari Kugemoto, Christian Muxica, Grusha Prasad, Brian Dillon, and Tal Linzen. 2024. [Large-scale benchmark yields no evidence that language model surprisal explains syntactic disambiguation difficulty](#). *Journal of Memory and Language*, 137:104510.
- Saketh Reddy Karra, Son The Nguyen, and Theja Tulabandhula. 2023. [Estimating the Personality of White-Box Language Models](#). *arXiv preprint*. ArXiv:2204.12000 [cs].
- Wolfgang Köhler. 1967. Gestalt psychology. *Psychologische forschung*, 31(1):XVIII–XXX.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo

- Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. [Solving Quantitative Reasoning Problems with Language Models](#). *arXiv preprint*. ArXiv:2206.14858 [cs].
- Kyle Mahowald, Evelina Fedorenko, Steven T Piantadosi, and Edward Gibson. 2013. Info/information theory: Speakers choose shorter words in predictive contexts. *Cognition*, 126(2):313–318.
- Christopher D. Manning, Kevin Clark, John Hewitt, Urvashi Khandelwal, and Omer Levy. 2020. [Emergent linguistic structure in artificial neural networks trained by self-supervision](#). *Proceedings of the National Academy of Sciences*, 117(48):30046–30054. Publisher: Proceedings of the National Academy of Sciences.
- James A. Michaelov and Benjamin K. Bergen. 2023. [Emergent inabilities? Inverse scaling over the course of pretraining](#). *arXiv preprint*. ArXiv:2305.14681 [cs].
- Marilù Miotto, Nicola Rossberg, and Bennett Kleinberg. 2022. [Who is GPT-3? An Exploration of Personality, Values and Demographics](#). *arXiv preprint*. ArXiv:2209.14338 [cs].
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altmenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, C. J. Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. [GPT-4 Technical Report](#). *arXiv preprint*. ArXiv:2303.08774 [cs].
- Qiwei Peng, Yekun Chai, and Xuhong Li. 2024. [HumanEval-XL: A Multilingual Code Generation Benchmark for Cross-lingual Natural Language Generalization](#). *arXiv preprint*. ArXiv:2402.16694 [cs].
- Martin J Pickering and Holly P Branigan. 1998. The representation of verbs: Evidence from syntactic priming in language production. *Journal of Memory and language*, 39(4):633–651.
- Shuofei Qiao, Yixin Ou, Ningyu Zhang, Xiang Chen, Yunzhi Yao, Shumin Deng, Chuanqi Tan, Fei Huang,

- and Huajun Chen. 2023. [Reasoning with Language Model Prompting: A Survey](#). *arXiv preprint*. ArXiv:2212.09597 [cs].
- Zhuang Qiu, Xufeng Duan, and Zhenguang Garry Cai. 2023. [Pragmatic Implicature Processing in ChatGPT](#).
- Qualtrics. 2024. [Qualtrics and all other qualtrics product or service names are registered trademarks or trademarks of qualtrics](#). Provo, UT, USA.
- Jennifer M Rodd, Belen Lopez Cutrin, Hannah Kirsch, Alessandra Millar, and Matthew H Davis. 2013. Long-term priming of the meanings of ambiguous words. *Journal of Memory and Language*, 68(2):180–198.
- Maarten Sap, Ronan LeBras, Daniel Fried, and Yejin Choi. 2023. [Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs](#). *arXiv preprint*. ArXiv:2210.13312 [cs].
- Iliia Shumailov, Zakhar Shumaylov, Yiren Zhao, Nicolas Papernot, Ross Anderson, and Yarin Gal. 2024. Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759.
- Murray Singer and Jackie Spear. 2015. Phantom recollection of bridging and elaborative inferences. *Discourse Processes*, 52(5-6):356–375.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. [Llama 2: Open Foundation and Fine-Tuned Chat Models](#). *arXiv preprint*. ArXiv:2307.09288 [cs].
- Sean Trott, Cameron Jones, Tyler Chang, James Michaelov, and Benjamin Bergen. 2023. [Do Large Language Models know what humans know?](#) *arXiv preprint*. ArXiv:2209.01515 [cs].
- Yuka Tsubota and Yoshinobu Kano. 2024. [Text Generation Indistinguishable from Target Person by Prompting Few Examples Using LLM](#). In *Proceedings of the 2nd International AIWolfDial Workshop*, pages 13–20, Tokyo, Japan. Association for Computational Linguistics.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. [HellaSwag: Can a Machine Really Finish Your Sentence?](#) *arXiv preprint*. ArXiv:1905.07830 [cs].
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, and Hongsheng Li. 2023. [Solving Challenging Math Word Problems Using GPT-4 Code Interpreter with Code-based Self-Verification](#). *arXiv preprint*. ArXiv:2308.07921 [cs] version: 1.

A Appendix

This section introduces the ten psycholinguistic experiments used to evaluate the humanlikeness of LLMs across multiple linguistic levels. Each experiment was designed to test a specific linguistic phenomenon and compare the performance of LLMs to human participants.

Sounds: sound-shape association People often associate specific sounds with certain shapes, a phenomenon known as sound symbolism. We tested whether LLMs, like humans, tend to link spiky-sounding words such as *takete* or *kiki* with spiky objects and round-sounding words like *maluma* or *bouba* with round objects.

Sounds: sound-gender association People can often guess if an unfamiliar name is male or female based on its sound. In English, women’s names more frequently end in vowels compared to men’s names. In this task, we asked participants to complete a preamble containing either a consonant-ending name (e.g., *Pelcrad* in 1a) or a vowel-ending novel name (e.g., *Pelcra* in 1b).

1a. Consonant-ending name: *Although Pelcrad was sick...*

1b. Vowel-ending name: *Although Pelcra was sick...*

Words: word length and predictivity Shorter words are suggested to make communication more efficient by carrying less information. If both humans and LLMs are sensitive to the relationship between word length and informativity, they

should prefer shorter words over longer ones with nearly identical meanings when completing sentence preambles that predicted the meaning of the word (making it less informative; e.g., 2a), compared to neutral sentence preambles (e.g., 2b)

2a. Predictive context: *Susan was very bad at algebra, so she hated... 1. math 2. mathematics*

2b. Neutral context: *Susan introduced herself to me as someone who loved... 1. math 2. mathematics*

Words: word meaning priming Many words have multiple meanings; for instance, *post* can refer to mail or a job. People update an ambiguous word's meaning based on recent exposure. We tested whether humans and LLMs similarly demonstrate word meaning priming phenomenon: Participants associated *post* with its job-related meaning more frequently after reading sentences using that context rather than synonyms' contexts (3a vs.3b).

3a. Word-meaning prime: *The man accepted the post in the accountancy firm.*

3b. Synonym prime: *The man accepted the job in the accountancy firm.*

Syntax: structural priming In structural priming, people tend to repeat syntactic structures they've recently encountered. We had participants complete prime preambles designed for either PO (prepositional-object dative structure, e.g., *The racing driver gave helpful mechanic wrench* to complete 4a) or DO (double-object dative structure, e.g., *The racing driver gave torn overall his mechanic* to complete 4b). Participants then completed target preamble which could be continued as either DO/PO. If structural priming is demonstrated, participants replicate structure of the prime preamble.

4a. DO-inducing prime preamble: *The racing driver showed the helpful mechanic ...*

4b. PO-inducing prime preamble: *The racing driver showed the torn overall ...*

4c. Target preamble: *The patient showed ...*

Syntax: syntactic ambiguity resolution The way people parse words into syntactic structures has garnered significant attention in psycholinguistics. For instance, in VP/NP ambiguity (e.g., *The ranger killed the poacher with the rifle*), people usually interpret the ambiguous prepositional phrase (PP, *with the rifle*) as modifying the verb phrase (VP, *killed the poacher*) rather than the noun phrase (NP, *the poacher*). However, contextual information can modulate this resolution: People are more likely to interpret ambiguous PPs as modifying NPs when

there are multiple possible referents (e.g., 5b) compared to when there is only a single referent (e.g., 5a). We examine how effectively LLMs use contextual information to resolve syntactic ambiguities and exhibit such modulation patterns.

5a. Single referent: *There was a hunter and a poacher. The hunter killed the dangerous poacher with a rifle not long after sunset. Who had a rifle, the hunter or the poacher?*

5b. Multiple referents: *There was a hunter and two poachers. The hunter killed the dangerous poacher with a rifle not long after sunset. Who had a rifle, the hunter or the poacher?*

Meaning: implausible sentence interpretation Listeners often need to recover intended messages from noise-corrupted input. Errors in production or comprehension can make a plausible sentence implausible by omitting (e.g., *to* omitted, 6a) or inserting words (e.g., *to* inserted, 6b). People may interpret an implausible sentence nonliterally if they believe it is noise-corrupted. who found that people more frequently reinterpret implausible DO sentences than PO sentences due to the likelihood of omissions over insertions. We tested whether people and LLMs similarly assume that implausible sentences result from noise corruption, with omissions being more likely than insertions.

6a. Implausible DO: *The mother gave the candle the daughter.*

6b. Implausible PO: *The mother gave the daughter to the candle.*

6c. Question: *Did the daughter receive something/someone?*

Meaning: semantic illusions People often overlook obvious errors in sentences. For instance, when asked (7a), many fail to notice that the question should refer to *Noah* instead of *Moses*. Such semantic illusions suggest that processing sentence meanings involves partial matches in semantic memory. We tested whether LLMs and people alike produce semantic illusions and are more likely to catch a weak imposter (e.g., *Adam*, less similar to *Noah*, 7b) than a strong imposter (e.g. *Morse*, more similar to *Noah*, 7a).

7a. Strong: *During the Biblical flood, how many animals of each kind did Moses take on the ark?*

7b. Weak: *During the Biblical flood, how many animals of each kind did Adam take on the ark?*

Discourse: implicit causality Certain verbs prompt people to associate causality with either the subject or the object within a sentence. For instance, stimulus-experiencer verbs like *scare* typ-

ically lead people to attribute causality to the subject (e.g., completing 8a as *Gary scared Anna because he was violent*), whereas experiencer-stimulus verbs like fear generally lead people to attribute causality to the object (e.g., completing 8b as *Gary feared Anna because she was violent*). We assessed whether LLMs, like humans, show similar patterns of causal attribution based on verb type.

8a. Stimulus-experiencer verb: *Gary scared Anna because...*

8b. Experiencer-stimulus verb: *Gary feared Anna because...*

Discourse: drawing inferences People make bridging inferences more frequently than elaborative inferences. Bridging inferences connect two pieces of information (after reading 9a, people infer that Sharon cut her foot) while elaborative inferences extrapolate from a single piece of information (people are less likely to make this inference after reading 9b). We examined how well an LLM aligns with human patterns of inference by comparing the bridging and elaborative conditions.

9a. Bridging: *While swimming in the shallow water near the rocks, Sharon stepped on a piece of glass. She called desperately for help, but there was no one around to hear her.*

9b. Elaborative: *While swimming in the shallow water near the rocks, Sharon stepped on a piece of glass. She had been looking for the watch that she misplaced while sitting on the rocks.*

Question: *Did she cut her foot?*